

Text Classification using Naive Bayes

Hiroshi Shimodaira*

11 February 2014

Text classification is the task of classifying documents by their content: that is, by the words of which they are comprised. Perhaps the best-known current text classification problem is email *spam filtering*: classifying email messages into spam and non-spam (ham).

1 Document models

Text classifiers often don't use any kind of deep representation about language: often a document is represented as a *bag of words*. (A bag is like a set that allows repeating elements.) This is an extremely simple representation: it only knows which words are included in the document (and how many times each word occurs), and throws away the word order!

Consider a document D , whose class is given by C . In the case of email spam filtering there are two classes $C = S$ (spam) and $C = H$ (ham). We classify D as the class which has the highest posterior probability $P(C|D)$, which can be re-expressed using Bayes' Theorem:

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)} \propto P(D|C)P(C). \quad (1)$$

We shall look at two probabilistic models of documents, both of which represent documents as a bag of words, using the Naive Bayes assumption. Both models represent documents using feature vectors whose components correspond to word types. If we have a vocabulary V , containing $|V|$ word types, then the feature vector dimension $d = |V|$.

Bernoulli document model: a document is represented by a feature vector with binary elements taking value 1 if the corresponding word is present in the document and 0 if the word is not present.

Multinomial document model: a document is represented by a feature vector with integer elements whose value is the frequency of that word in the document.

Example: Consider the vocabulary:

$$V = \{\text{blue, red, dog, cat, biscuit, apple}\}.$$

In this case $|V| = d = 6$. Now consider the (short) document "the blue dog ate a blue biscuit". If \mathbf{d}^B is the Bernoulli feature vector for this document, and \mathbf{d}^M is the multinomial feature vector, then we

*Heavily based on notes inherited from Steve Renals and Iain Murray.

would have:

$$\mathbf{d}^B = (1, 0, 1, 0, 1, 0)^T$$

$$\mathbf{d}^M = (2, 0, 1, 0, 1, 0)^T$$

To classify a document we use equation (1), which requires estimating the likelihoods of the document given the class, $P(D|C)$ and the class prior probabilities $P(C)$. To estimate the likelihood, $P(D|C)$, we use the Naive Bayes assumption applied to whichever of the two document models we are using.

2 The Bernoulli document model

As mentioned above, in the Bernoulli model a document is represented by a binary vector, which represents a point in the space of words. If we have a vocabulary V containing a set of $|V|$ words, then the i th dimension of a document vector corresponds to word w_i in the vocabulary. Let \mathbf{b}_i be the feature vector for the i th document D^i ; then the i th element of \mathbf{b}_i , written b_{ii} , is either 0 or 1 representing the absence or presence of word w_i in the i th document.

Let $P(w_i|C)$ be the probability of word w_i occurring in a document of class C ; the probability of w_i not occurring in a document of this class is given by $(1 - P(w_i|C))$. If we make the naive Bayes assumption, that the probability of each word occurring in the document is independent of the occurrences of the other words, then we can write the document likelihood $P(D^i|C)$ in terms of the individual word likelihoods $P(w_i|C)$:

$$P(D^i|C) \sim P(\mathbf{b}_i|C) = \prod_{i=1}^{|V|} [b_{ii}P(w_i|C) + (1 - b_{ii})(1 - P(w_i|C))]. \quad (2)$$

This product goes over all words in the vocabulary. If word w_i is present, then $b_{ii} = 1$ and the required probability is $P(w_i|C)$; if word w_i is not present, then $b_{ii} = 0$ and the required probability is $1 - P(w_i|C)$. We can imagine this as a model for generating document feature vectors of class C , in which the document feature vector is modelled as a collection of $|V|$ weighted coin tosses, the i th having a probability of success equal to $P(w_i|C)$.

The *parameters* of the likelihoods are the probabilities of each word given the document class $P(w_i|C)$; the model is also parameterised by the prior probabilities, $P(C)$. We can learn (estimate) these parameters from a training set of documents labelled with class $C = k$. Let $n_k(w_i)$ be the number of documents of class $C = k$ in which w_i is observed; and let N_k be the total number of documents of that class. Then we can estimate the parameters of the word likelihoods as,

$$\hat{P}(w_i | C = k) = \frac{n_k(w_i)}{N_k}, \quad (3)$$

the relative frequency of documents of class $C = k$ that contain word w_i . If there are N documents in total in the training set, then the prior probability of class $C = k$ may be estimated as the relative frequency of documents of class $C = k$:

$$\hat{P}(C = k) = \frac{N_k}{N}. \quad (4)$$

Thus given a training set of documents (each labelled with a class), and a set of K classes, we can estimate a Bernoulli text classification model as follows:

1. Define the vocabulary V ; the number of words in the vocabulary defines the dimension of the feature vectors
2. Count the following in the training set:
 - N the total number of documents
 - N_k the number of documents labelled with class $C=k$, for $k=1, \dots, K$
 - $n_k(w_i)$ the number of documents of class $C=k$ containing word w_i for every class and for each word in the vocabulary
3. Estimate the likelihoods $P(w_i | C=k)$ using equation (3)
4. Estimate the priors $P(C=k)$ using equation (4)

To classify an unlabelled document D^j , we estimate the posterior probability for each class combining equations (1) and (2):

$$\begin{aligned}
 P(C|D^j) &= P(C|\mathbf{b}_j) \\
 &\propto P(\mathbf{b}_j|C)P(C) \\
 &\propto P(C) \prod_{i=1}^{|V|} [b_{ji}P(w_i|C) + (1-b_{ji})(1-P(w_i|C))].
 \end{aligned} \quad (5)$$

Example

Consider a set of documents, each of which is related either to *Sports* (S) or to *Informatics* (I). Given a training set of 11 documents, we would like to estimate a Naive Bayes classifier, using the Bernoulli document model, to classify unlabelled documents as S or I .

We define a vocabulary of eight words:

$$V = \begin{bmatrix} w_1 = \text{goal,} \\ w_2 = \text{tutor,} \\ w_3 = \text{variance,} \\ w_4 = \text{speed,} \\ w_5 = \text{drink,} \\ w_6 = \text{defence,} \\ w_7 = \text{performance,} \\ w_8 = \text{field} \end{bmatrix}$$

Thus each document is represented as an 8-dimensional binary vector.

The training data is presented below as a matrix for each class, in which each row represents an 8-dimensional document vector

$$\mathbf{B}^{\text{Sport}} = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$\mathbf{B}^{\text{Inf}} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

Classify the following into Sports or Informatics using a Naive Bayes classifier.

$$1. \mathbf{b}_1 = (1, 0, 0, 1, 1, 1, 0, 1)^T$$

$$2. \mathbf{b}_2 = (0, 1, 1, 0, 1, 0, 1, 0)^T$$

Solution:

The total number of documents in the training set $N=11$; $N_S=6$, $N_I=5$

Using (4), we can estimate the prior probabilities from the training data as:

$$P(S) = \frac{6}{11}; \quad P(I) = \frac{5}{11}$$

The word counts in the training data are:

$$\begin{array}{ll} n_S(w_1) = 3 & n_S(w_2) = 1 \\ n_S(w_3) = 2 & n_S(w_4) = 3 \\ n_S(w_5) = 3 & n_S(w_6) = 4 \\ n_S(w_7) = 4 & n_S(w_8) = 4 \end{array}$$

$$\begin{array}{ll} n_I(w_1) = 1 & n_I(w_2) = 3 \\ n_I(w_3) = 3 & n_I(w_4) = 1 \\ n_I(w_5) = 1 & n_S(w_6) = 1 \\ n_S(w_7) = 3 & n_S(w_8) = 1 \end{array}$$

The we can estimate the word likelihoods using (3)

$$\begin{array}{ll} P(w_1|S) = \frac{1}{2} & P(w_2|S) = \frac{1}{6} \\ P(w_3|S) = \frac{1}{3} & P(w_4|S) = \frac{1}{2} \\ P(w_5|S) = \frac{1}{2} & P(w_6|S) = \frac{2}{3} \\ P(w_7|S) = \frac{2}{3} & P(w_8|S) = \frac{2}{3} \end{array}$$

And for class I :

$$\begin{aligned} P(w_1|I) &= \frac{1}{5} & P(w_2|I) &= \frac{3}{5} \\ P(w_3|I) &= \frac{3}{5} & P(w_4|I) &= \frac{1}{5} \\ P(w_5|I) &= \frac{1}{5} & P(w_6|I) &= \frac{1}{5} \\ P(w_7|I) &= \frac{3}{5} & P(w_8|I) &= \frac{1}{5} \end{aligned}$$

We use (5) to compute the posterior probabilities of the two test vectors and hence classify them.

$$1. \mathbf{b}_1 = (1, 0, 0, 1, 1, 1, 0, 1)^T$$

$$\begin{aligned} P(S|\mathbf{b}_1) &\propto P(S) \prod_{t=1}^8 [b_{1t}P(w_t|S) + (1 - b_{1t})(1 - P(w_t|S))] \\ &\propto \frac{6}{11} \left(\frac{1}{2} \times \frac{5}{6} \times \frac{2}{3} \times \frac{1}{2} \times \frac{1}{2} \times \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} \right) = \frac{5}{891} = 5.6 \times 10^{-3} \\ P(I|\mathbf{b}_1) &\propto P(I) \prod_{t=1}^8 [b_{1t}P(w_t|I) + (1 - b_{1t})(1 - P(w_t|I))] \\ &\propto \frac{5}{11} \left(\frac{1}{5} \times \frac{2}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{1}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{1}{5} \right) = \frac{8}{859375} = 9.3 \times 10^{-6} \end{aligned}$$

Classify this document as S .

$$2. \mathbf{b}_2 = (0, 1, 1, 0, 1, 0, 1, 0)^T$$

$$\begin{aligned} P(S|\mathbf{b}_2) &\propto P(S) \prod_{t=1}^8 [b_{2t}P(w_t|S) + (1 - b_{2t})(1 - P(w_t|S))] \\ &\propto \frac{6}{11} \left(\frac{1}{2} \times \frac{1}{6} \times \frac{1}{3} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} \right) = \frac{12}{14256} = 8.4 \times 10^{-4} \\ P(I|\mathbf{b}_2) &\propto P(I) \prod_{t=1}^8 [b_{2t}P(w_t|I) + (1 - b_{2t})(1 - P(w_t|I))] \\ &\propto \frac{5}{11} \left(\frac{4}{5} \times \frac{3}{5} \times \frac{3}{5} \times \frac{4}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{4}{5} \right) = \frac{34560}{4296875} = 8.0 \times 10^{-3} \end{aligned}$$

Classify as I .

3 The multinomial distribution

Before discussing the multinomial document model, it is important to be familiar with the multinomial distribution.

We first need to be able to count the number of distinct arrangements of a set of items, when some of the items are *indistinguishable*. For example: Using all the letters, how many distinct sequences

can you make from the word “Mississippi”? There are 11 letters to permute, but “i” and “s” occur four times and “p” twice. If these letters were distinct (e.g., if they were labelled i_1, i_2 , etc.) then there would be $11!$ permutations. However of these permutations there are $4!$ that are the same if the subscripts are removed from the “i”s. This means that we can reduce the size of the total sample space by a factor of $4!$ to take account of the four occurrences of “i”. Likewise there is a factor of $4!$ for “s” and a factor of $2!$ for “p” (and a factor of $1!$ for “m”). This gives the total number distinct permutations as:

$$\frac{11!}{4! 4! 2! 1!} = 34650$$

Generally if we have n items of d types, with n_1 of type 1, n_2 of type 2 and n_d of type d (such that $n_1 + n_2 + \dots + n_d = n$), then the number of distinct permutations is given by:

$$\frac{n!}{n_1! n_2! \dots n_d!}$$

These numbers are called the *multinomial coefficients*.

Now suppose a population contains items of $d \geq 2$ different types and that the proportion of items that are of type t is p_t ($t = 1, \dots, d$), with

$$\sum_{t=1}^d p_t = 1 \quad p_t > 0, \text{ for all } t.$$

Suppose n items are drawn at random (with replacement) and let x_t denote the number of items of type t . The vector $\mathbf{x} = (x_1, \dots, x_d)^T$ has a *multinomial distribution* with parameters n and p_1, \dots, p_d , defined by:

$$\begin{aligned} P(\mathbf{x}) &= \frac{n!}{x_1! x_2! \dots x_d!} p_1^{x_1} p_2^{x_2} \dots p_d^{x_d} \\ &= \frac{n!}{\prod_{t=1}^d x_t!} \prod_{t=1}^d p_t^{x_t} \end{aligned} \quad (6)$$

The $\prod_{t=1}^d p_t^{x_t}$ product gives the probability of one sequence of outcomes with counts \mathbf{x} . The multinomial coefficient, counts the number of such sequences that there are.

4 The multinomial document model

In the multinomial document model, the document feature vectors capture the frequency of words, not just their presence or absence. Let \mathbf{x}_i be the multinomial model feature vector for the i th document D^i . The t th element of \mathbf{x}_i , written x_{it} , is the count of the number of times word w_t occurs in document D^i . Let $n_i = \sum_t x_{it}$ be the total number of words in document D^i .

Let $P(w_t|C)$ again be the probability of word w_t occurring in class C , this time estimated using the word frequency information from the document feature vectors. We again make the naive Bayes assumption, that the probability of each word occurring in the document is independent of the occurrences of the other words. We can then write the document likelihood $P(D^i|C)$ as a multinomial distribution (equation 6), where the number of draws corresponds to the length of the document, and the proportion

of drawing item t is the probability of word type t occurring in a document of class C , $P(w_t|C)$.

$$P(D^j|C) \sim P(\mathbf{x}_j|C) = \frac{n_j!}{\prod_{t=1}^{|V|} x_{jt}!} \prod_{t=1}^{|V|} P(w_t|C)^{x_{jt}} \\ \propto \prod_{t=1}^{|V|} P(w_t|C)^{x_{jt}}. \quad (7)$$

We often won't need the normalisation term ($n_j! / \prod_t x_{jt}!$), because it does not depend on the class, C . The numerator of the right hand side of this expression can be interpreted as the product of word likelihoods for each word in the document, with repeated words taking part for each repetition.

As for the Bernoulli model, the parameters of the likelihood are the probabilities of each word given the document class $P(w_t|C)$, and the model parameters also include the prior probabilities $P(C)$. To estimate these parameters from a training set of documents labelled with class $C = k$, let z_{ik} be an indicator variable which equals 1 when D^i has class $C = k$, and equals 0 otherwise. If N is again the total number of documents, then we have:

$$\hat{P}(w_t | C=k) = \frac{\sum_{i=1}^N x_{it} z_{ik}}{\sum_{s=1}^{|V|} \sum_{i=1}^N x_{is} z_{ik}}, \quad (8)$$

an estimate of the probability $P(w_t | C=k)$ as the relative frequency of w_t in documents of class $C = k$ with respect to the total number of words in documents of that class.

The prior probability of class $C = k$ is estimated as before (equation 4).

Thus given a training set of documents (each labelled with a class) and a set of K classes, we can estimate a multinomial text classification model as follows:

1. Define the vocabulary V ; the number of words in the vocabulary defines the dimension of the feature vectors.
2. Count the following in the training set:
 - N the total number of documents,
 - N_k the number of documents labelled with class $C = k$, for each class $k = 1, \dots, K$,
 - x_{it} the frequency of word w_t in document D^i , computed for every word w_t in V .
3. Estimate the likelihoods $P(w_t | C=k)$ using (8).
4. Estimate the priors $P(C=k)$ using (4).

To classify an unlabelled document D^j , we estimate the posterior probability for each class combining (1) and (7):

$$P(C|D^j) = P(C|\mathbf{x}_j) \\ \propto P(\mathbf{x}_j|C) P(C) \\ \propto P(C) \prod_{t=1}^{|V|} P(w_t|C)^{x_{jt}}. \quad (9)$$

Unlike the Bernoulli model, words that do not occur in the document (i.e., for which $x_{it} = 0$) do not affect the probability (since $p^0 = 1$). Thus we can write the posterior probability in terms of words u which occur in the document:

$$P(C|D^j) \propto P(C) \prod_{h=1}^{\text{len}(D^j)} P(u_h|C)$$

Where u_h is the h th word in document D^j .

5 The Zero Probability Problem

A drawback of relative frequency estimates—equation (8) for the multinomial model—is that zero counts result in estimates of zero probability. This is a bad thing because the Naive Bayes equation for the likelihood (7) involves taking a product of probabilities: if any one of the terms of the product is zero, then the whole product is zero. This means that the probability of the document belonging to the class in question is zero—which means it is impossible.

Just because a word does not occur in a document class in the training data does not mean that it cannot occur in any document of that class.

The problem is that equation (8) *underestimates* the likelihoods of words that do not occur in the data. Even if word w is not observed for class $C = k$ in the training set, we would still like $P(w | C=k) > 0$. Since probabilities must sum to 1, if unobserved words have underestimated probabilities, then those words that are observed must have overestimated probabilities. Therefore, one way to alleviate the problem is to remove a small amount of probability allocated to observed events and distribute this across the unobserved events. A simple way to do this, sometimes called *Laplace's law of succession* or *add one smoothing*, adds a count of one to each word type. If there are W word types in total, then equation (8) may be replaced with:

$$P_{\text{Lap}}(w_t | C=k) = \frac{1 + \sum_{i=1}^N x_{it} z_{ik}}{|V| + \sum_{s=1}^{|V|} \sum_{i=1}^N x_{is} z_{ik}} \quad (10)$$

The denominator was increased to take account of the $|V|$ extra “observations” arising from the “add 1” term, ensuring that the probabilities are still normalised.

Question: The Bernoulli document model can also suffer from the zero probability model. How would you apply add one smoothing in this case?

6 Comparing the two models

The Bernoulli and the multinomial document models are both based on a bag of words. However there are a number of differences, which we summarise here:

1. Underlying model of text:

Bernoulli: a document can be thought of as being generated from a multidimensional Bernoulli distribution: the probability of a word being present can be thought of as a (weighted) coin flip with probability $P(w_t|C)$.

Multinomial: a document is formed by drawing words from a multinomial distribution: you can think of obtaining the next word in the document by rolling a (weighted) $|V|$ -sided dice with probabilities $P(w_t|C)$.

2. Document representation:

Bernoulli: binary vector, elements indicating presence or absence of a word.

Multinomial: integer vector, elements indicating frequency of occurrence of a word.

3. Multiple occurrences of words:

Bernoulli: ignored.

Multinomial: taken into account.

4. Behaviour with document length:

Bernoulli: best for short documents.

Multinomial: longer documents are OK.

5. Behaviour with “the”:

Bernoulli: since “the” is present in almost every document, $P(\text{“the”}|C) \sim 1.0$.

Multinomial: since probabilities are based on relative frequencies of word occurrence in a class, $P(\text{“the”}|C) \sim 0.05$.

6. Non-occurring words:

Bernoulli: affect the document probabilities.

Multinomial: do not affect the document probabilities.

7 Conclusion

In this chapter we have shown how the Naive Bayes approximation can be used for document classification, by constructing distributions over words. The classifiers require a *document model* to estimate $P(\text{document} | \text{class})$. We looked at two document models that we can use with the Naive Bayes approximation:

- **Bernoulli document model:** a document is represented by a binary feature vector, whose elements indicate absence or presence of corresponding word in the document.
- **Multinomial document model:** a document is represented by an integer feature vector, whose elements indicate frequency of corresponding word in the document.